



# A High Throughput Polynomial and Rational Function Approximations Evaluator

Nicolas Brisebarre, George Constantinides, Miloš Ercegovac, Silviu-Ioan Filip, Matei Istoan, Jean-Michel Muller

## ► To cite this version:

Nicolas Brisebarre, George Constantinides, Miloš Ercegovac, Silviu-Ioan Filip, Matei Istoan, et al.. A High Throughput Polynomial and Rational Function Approximations Evaluator. ARITH 2018 - 25th IEEE Symposium on Computer Arithmetic, Jun 2018, Amherst, MA, United States. pp.99-106, 10.1109/ARITH.2018.8464778 . hal-01774364

**HAL Id: hal-01774364**

**<https://inria.hal.science/hal-01774364>**

Submitted on 23 Apr 2018

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# A High Throughput Polynomial and Rational Function Approximations Evaluator

Nicolas Brisebarre<sup>1</sup>, George Constantinides<sup>2</sup>, Miloš Ercegovic<sup>3</sup>, Silviu-Ioan Filip<sup>4</sup>, Matei Istoan<sup>2</sup>, and Jean-Michel Muller<sup>1</sup>

<sup>1</sup>Laboratoire LIP, CNRS, ENS Lyon, INRIA, Univ. Claude Bernard Lyon 1, Lyon, France, f-name.l-name@ens-lyon.fr

<sup>2</sup>Electrical and Electronic Engineering, Imperial College London, SW7 2AZ, United Kingdom, {g.constantinides,m.istoan}@imperial.ac.uk

<sup>3</sup>Computer Science Department, UCLA, Los Angeles, CA, 90095, USA, milos@cs.ucla.edu

<sup>4</sup>Univ Rennes, Inria, CNRS, IRISA, F-35000 Rennes, silviu.filip@inria.fr

April 23, 2018

## Abstract

We present an automatic method for the evaluation of functions via polynomial or rational approximations and its hardware implementation, on FPGAs. These approximations are evaluated using Ercegovic's iterative E-method adapted for FPGA implementation. The polynomial and rational function coefficients are optimized such that they satisfy the constraints of the E-method. We present several examples of practical interest; in each case a resource-efficient approximation is proposed and comparisons are made with alternative approaches.

## 1 Introduction

We aim at designing a system able to approximate (in software) and then evaluate (in hardware) any regular-enough function. More precisely, we try to minimize the sup norm of the difference between the function and the approximation in a given interval.

For particular functions, *ad hoc* solutions such as CORDIC [1] or some specific tabulate-and-compute algorithms [2] can be used. For low precision cases, table-based methods [3–5] methods are of interest. However, in the general case, piecewise approximations by polynomial or rational functions are the only reasonable solution. From a theoretical point of view, rational functions are very attractive, mainly because they can reproduce function behaviors (such as asymptotes, finite limits at  $\pm\infty$ ) that polynomials do not satisfy. However, for software implementation, polynomials are frequently preferred to rational functions, because the latency of division is larger than the latency of multiplication. We aim at checking if rational approximations are of interest in hardware implementations. To help in the comparison of polynomial and rational approximations in hardware we use an algorithm, due to Ercegovic [6, 7], called the E-method, that makes it possible to evaluate a degree- $n$  polynomial, or a rational function of degree- $n$  numerator and denominator at a similar cost without requiring division.

The E-method solves diagonally-dominant linear systems using a left-to-right digit-by-digit approach and has a simple and regular hardware implementation. It maps the problem of evaluating a polynomial or rational function into a linear system. The linear system corresponding to a given function does not necessarily satisfy the conditions of diagonal dominance. For polynomials, changes of variables allow one to satisfy the conditions. This is not the case for rational functions. There is however a family of rational functions, called E-fractions, that can be evaluated with the E-method in time proportional to the desired precision. One of our aims is, given a function, to decide whether it is better to approximate it by a polynomial or by an E-fraction. Furthermore, we want to design approximations whose coefficients satisfy some constraints (such as being exactly representable in a given format). We introduce algorithmic improvements with respect to [8] for computing E-fractions. We

present a circuit generator for the E-method and compare its implementation on an FPGA with FloPoCo polynomial designs [9] for several examples of practical interest. Since FloPoCo designs are pipelined (unrolled), we focus on an unrolled design of the E-method.

## 1.1 An Overview of the E-method

The E-method evaluates a polynomial  $P_\mu(x)$  or a rational function  $R_{\mu,\nu}(x)$  by mapping it into a linear system. The system is solved using a left-to-right digit-by-digit approach, in a radix  $r$  representation system, on a regular hardware. For a result of  $m$  digits, in the range  $(-1, 1)$ , the computation takes  $m$  iterations. The first component of the solution vector corresponds to the value of  $P_\mu(x)$  or  $R_{\mu,\nu}(x)$ . Let

$$R_{\mu,\nu}(x) = \frac{P_\mu(x)}{Q_\nu(x)} = \frac{p_\mu x^\mu + p_{\mu-1} x^{\mu-1} + \dots + p_0}{q_\nu x^\nu + q_{\nu-1} x^{\nu-1} + \dots + q_1 x + 1}$$

where the  $p_i$ 's and  $q_i$ 's are real numbers. Let  $n = \max\{\mu, \nu\}$ ,  $p_j = 0$  for  $\mu + 1 \leq j \leq n$ , and  $q_j = 0$  for  $\nu + 1 \leq j \leq n$ . According to the E-method  $R_{\mu,\nu}(x)$  is mapped to a linear system  $L : \mathbf{A} \times \mathbf{y} = \mathbf{b}$ :

$$\begin{bmatrix} 1 & -x & 0 & \dots & 0 \\ q_1 & 1 & -x & 0 & \dots & 0 \\ q_2 & 0 & 1 & -x & \dots & 0 \\ & & \ddots & \ddots & & \vdots \\ \vdots & & & & \ddots & 0 \\ q_{n-1} & & & & 1 & -x \\ q_n & & \dots & 0 & 1 \end{bmatrix} \begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ \vdots \\ \vdots \\ y_{n-1} \\ y_n \end{bmatrix} = \begin{bmatrix} p_0 \\ p_1 \\ p_2 \\ \vdots \\ \vdots \\ p_{n-1} \\ p_n \end{bmatrix} \quad (1)$$

so that  $y_0 = R_{\mu,\nu}(x)$ . Likewise,  $y_0 = P_\mu(x)$  when all  $q_i = 0$ .

The components of the solution vector  $\mathbf{y} = [y_0, y_1, \dots, y_n]^t$  are computed, digit-by-digit, the most-significant digit first, by means of the following vector iteration:

$$\mathbf{w}^{(j)} = r \times \left[ \mathbf{w}^{(j-1)} - \mathbf{A} \mathbf{d}^{(j-1)} \right], \quad (2)$$

for  $j = 1, \dots, m$ , where  $m$  is the desired precision of the result. The term  $\mathbf{w}^{(j)}$  is the vector residual in iteration  $j$  with  $\mathbf{w}^{(0)} = [p_0, p_1, \dots, p_n]^t$ . The solution  $\mathbf{y}$  is produced as a sequence of digit vectors:  $\mathbf{d}^{(j-1)} = [d_1^{(j-1)}, \dots, d_n^{(j-1)}]^t$  – a digit vector obtained in iteration  $j-1$  and used in iteration  $j$ . After  $m$  iterations,  $y_k = \sum_{j=1}^m d_k^{(j)} r^{-j}$ . The digits of the solution components  $y_0, y_1, \dots, y_n$  are computed using very simple scalar recurrences. Note that all multiplications in these recurrences use  $m \times 1$  multipliers and that division required by the rational function is not explicitly performed.

$$w_i^{(j)} = r \times \left[ w_i^{(j-1)} - q_i d_0^{(j-1)} - d_i^{(j-1)} + d_{i+1}^{(j-1)} x \right], \quad (3)$$

$$w_0^{(j)} = r \times \left[ w_0^{(j-1)} - d_0^{(j-1)} + d_1^{(j-1)} x \right] \quad (4)$$

and

$$w_n^{(j)} = r \times \left[ w_n^{(j-1)} - d_n^{(j-1)} - q_n d_0^{(j-1)} \right]. \quad (5)$$

Initially,  $\mathbf{d}^{(0)} = \mathbf{0}$ . The radix- $r$  digits  $d_i^{(j)}$  are in the redundant signed digit-set  $D_\rho = \{-\rho, \dots, 0, 1, \dots, \rho\}$  with  $r/2 \leq \rho \leq r-1$ . If  $\rho = r/2$ ,  $D_\rho$  is called *minimally redundant*, and if  $\rho = r-1$ , it is *maximally redundant*. The choice of redundancy is determined by design considerations. The radix of computation is  $r = 2^k$  so that internally radix-2 arithmetic is used. The residuals, in general, are in a redundant form to reduce the iteration time. Since the target is an FPGA technology which provides fast carry chains, we have non-redundant residuals. The digits  $d_i^{(j)}$  are selected so that the residuals  $|w_i^{(j)}|$  remain bounded. The digit selection is performed by rounding the residuals  $w_i^{(j)}$  to a single signed digit, following [7, 10]:

$$d_i^{(j)} = S(w_i^{(j)}) = \begin{cases} \text{sign}(w_i^{(j)}) \times \left\lfloor \left| \frac{w_i^{(j)}}{r} \right| + \frac{1}{2} \right\rfloor, & \text{if } |w_i^{(j)}| \leq \rho, \\ \text{sign}(w_i^{(j)}) \times \left\lfloor \left| \frac{w_i^{(j)}}{r} \right| \right\rfloor, & \text{otherwise.} \end{cases}$$

The selection is performed using a low-precision estimate  $\widehat{w}_i^{(j)}$  of  $w_i^{(j)}$ , obtained by truncating  $w_i^{(j)}$  to one fractional bit.

Since the matrices considered here have 1s on the diagonal, a necessary condition for convergence is  $\sum_{j \neq i} |a_{i,j}| < 1$ . Specifically,

$$\begin{cases} \forall i, |p_i| \leq \xi, \\ \forall i, |x| + |q_i| \leq \alpha, \\ |w_i^{(j)} - \widehat{w}_i^{(j)}| \leq \Delta/2. \end{cases} \quad (6)$$

where the bounds  $\xi$ ,  $\alpha$ , and  $\Delta$  satisfy [7]:

$$\xi = \frac{1}{2}(1 + \Delta), \quad 0 < \Delta < 1, \quad \alpha \leq (1 - \Delta)/(2r) \quad (7)$$

for maximally redundant digit sets used here. While the constraints (7) may seem restrictive, for polynomials, scaling techniques make it possible to satisfy them. However, this is not the case for all rational functions. To remove this limitation the authors of [8] have suggested the derivation of rational functions, called *simple E-fractions*, which are products of a power of 2 by a fraction that satisfies (7). In this work we make further improvements to the rational functions based on E-fractions.

## 1.2 Outline of the paper

In Section 2, we discuss the effective generation of simple E-fractions, whose coefficients are exactly representable in a given format. Section 3 presents a hardware implementation of the E-method that targets FPGAs. In Section 4 we present and discuss some examples in various situations. We also present a comparison with FloPoCo implementations.

# 2 Effective computation of simple E-fractions

We show how to compute a simple E-fraction with fixed-point or floating-point coefficients. A first step (see Section 2.1), yields a simple E-fraction approximation with real coefficients to a function  $f$ . In [8], linear programming (LP) is used. Here, we use faster tools from approximation theory. This allows us to quickly check how far the approximation error of this E-fraction is from the optimal error of the minimax approximation (obtained using the Remez algorithm [11, 12]), and how far it is from the error that an E-polynomial, with the same implementation cost, can yield. If this comparison suggests that it is more advantageous to work with an E-fraction, we use the Euclidean lattice basis reduction approach from [8] for computing E-fractions with machine-number coefficients. We introduce in Section 2.2.2 a trick that improves its output.

## 2.1 Real approximation step

Let  $f$  be a continuous function defined on  $[a, b]$ . Let  $\mu, \nu \in \mathbb{N}$  be given and let  $\mathbb{R}_{\mu, \nu}(x) = \{P/Q : P = \sum_{k=0}^{\mu} p_k x^k, Q = \sum_{k=0}^{\nu} q_k x^k, p_0, \dots, p_{\mu}, q_0, \dots, q_{\nu} \in \mathbb{R}\}$ . The aim is to compute a good rational fraction approximant  $R \in \mathbb{R}_{\mu, \nu}(x)$ , with respect to the supremum norm defined by  $\|g\| = \sup_{x \in [a, b]} |g(x)|$ , to  $f$  such that the real coefficients of  $R$  (or  $R$  divided by some fixed power of 2) satisfy the constraints imposed by the E-method.

As done in [8], we can first apply the rational version of the Remez exchange algorithm [11, p. 173] to get  $R^*$ , the best possible rational approximant to  $f$  among the elements of  $\mathbb{R}_{\mu, \nu}(x)$ . This algorithm can fail if  $R^*$  is degenerate or the choice of starting nodes is not good enough.

To bypass these issues, we develop the following process. It can be viewed as a Remez-like method of the *first type*, following ideas described in [11, p. 96–97] and [13]. It directly computes best real coefficient E-fractions with magnitude constraints on the denominator coefficients. If we remove these constraints, it will compute the minimax rational approximation, even when the Remez exchange algorithm fails.

We first show how to solve the problem over  $X$ , a finite discretization of  $[a, b]$ . We apply a modified version (with denominator coefficient magnitude constraints) of the differential correction (DC) algorithm introduced in [14]. It is given by Algorithm 1. System (8) is an LP problem and can be solved in practice very efficiently using a simplex-based LP solver. Convergence of this **EDiffCorr** procedure can be shown using an identical argument to the convergence proofs of the original DC algorithm [15, 16] (see Appendix A).

---

**Algorithm 1** E-fraction EDiffCorr algorithm

---

**Input:**  $f \in \mathcal{C}([a, b])$ ,  $\mu, \nu \in \mathbb{N}$ , finite set  $X \subseteq [a, b]$  with  $|X| > \mu + \nu$ , threshold  $\varepsilon > 0$ , coefficient magnitude bound  $d > 0$

**Output:** approximation  $R(x) = \frac{\sum_{k=0}^{\mu} p_k x^k}{1 + \sum_{k=1}^{\nu} q_k x^k}$  of  $f$  over  $X$  s.t.  $\max_{1 \leq k \leq \nu} |q_k| \leq d$

// Initialize the iterative procedure ( $R = P/Q$ )

1:  $R \leftarrow 1$

2: **repeat**

3:  $\delta \leftarrow \max_{x \in X} |f(x) - R(x)|$

4: find  $R_{\text{new}} = P_{\text{new}}/Q_{\text{new}} = \frac{\sum_{k=0}^{\mu} p'_k x^k}{1 + \sum_{k=1}^{\nu} q'_k x^k}$  such that the expression

$$\max_{x \in X} \left\{ \frac{|f(x)Q_{\text{new}}(x) - P_{\text{new}}(x)| - \delta Q_{\text{new}}(x)}{Q(x)} \right\} \quad (8)$$

subject to  $\max_{1 \leq k \leq \nu} |q'_k| \leq d$ , is minimized

5:  $\delta_{\text{new}} \leftarrow \max_{x \in X} |f(x) - R_{\text{new}}(x)|$

6:  $R \leftarrow R_{\text{new}}$

7: **until**  $|\delta - \delta_{\text{new}}| < \varepsilon$

---

**Algorithm 2** E-fraction Remez algorithm

---

**Input:**  $f \in \mathcal{C}([a, b])$ ,  $\mu, \nu \in \mathbb{N}$ , finite set  $X \subseteq [a, b]$  with  $|X| > \mu + \nu$ , threshold  $\varepsilon > 0$ , coefficient magnitude bound  $d > 0$

**Output:** approximation  $R^*(x) = \frac{\sum_{k=0}^{\mu} p_k^* x^k}{1 + \sum_{k=1}^{\nu} q_k^* x^k}$  of  $f$  over  $[a, b]$  s.t.  $\max_{1 \leq k \leq \nu} |q_k^*| \leq d$

// Compute best E-fraction approximation over  $X$  using a

// modified version of the differential correction algorithm

1:  $R^* \leftarrow \text{EDiffCorr}(f, \mu, \nu, X, \varepsilon, d)$

2:  $\delta^* \leftarrow \max_{x \in X} |f(x) - R^*(x)|$

3:  $\Delta^* \leftarrow \max_{x \in [a, b]} |f(x) - R^*(x)|$

4: **while**  $\Delta^* - \delta^* > \varepsilon$  **do**

5:  $x_{\text{new}} \leftarrow \operatorname{argmax}_{x \in [a, b]} |f(x) - R^*(x)|$

6:  $X \leftarrow X \cup \{x_{\text{new}}\}$

7:  $R^* \leftarrow \text{EDiffCorr}(f, \mu, \nu, X, \varepsilon, d)$

8:  $\delta^* \leftarrow \max_{x \in X} |f(x) - R^*(x)|$

9:  $\Delta^* \leftarrow \max_{x \in [a, b]} |f(x) - R^*(x)|$

10: **end while**

---

To address the problem over  $[a, b]$ , Algorithm 2 solves a series of best E-fraction approximation problems on a discrete subset  $X$  of  $[a, b]$ , where  $X$  increases at each iteration by adding a point where the current residual term achieves its global maximum.

Our current experiments suggest that the speed of convergence for Algorithm 2 is linear. We can potentially decrease the number of iterations by adding to  $X$  more local extrema of the residual term at each iteration. Other than its speed compared to the LP approach from [8], Algorithm 2 will generally converge to the best E-fraction approximation with real coefficients over  $[a, b]$ , and not on a discretization of  $[a, b]$ .

Once  $R^*$  is computed, we determine the least integer  $s$  such that the coefficients of the numerator of  $R^*$  divided by  $2^s$  fulfill the first condition of (6). It gives us a decomposition  $R^*(x) = 2^s R_s(x)$ .  $R_s$  is thus a rescaled version of  $R$ . We take  $f_s = 2^{-s} f$  to be the corresponding rescaling of  $f$ . The magnitude bound  $d$  is usually equal to  $\alpha - \max(|a|, |b|)$ , allowing the denominator coefficients to be valid with respect to the second constraint of (6).

Both Algorithm 1 and 2 can be modified to compute weighted error approximations, that is, work with a norm of the form  $\|g\| = \max_{x \in [a, b]} |w(x)g(x)|$ , where  $w$  is a continuous and positive weight function over  $[a, b]$ . This is useful, for instance, when targeting relative error approximations. The changes are minimal and consist only of introducing the weight factor in the error computations in lines 3, 5 of Algorithm 1, lines 2, 3, 5, 8, 9 of

Algorithm 2 and changing (8) with

$$\max_{x \in X} \left\{ \frac{w(x) |f(x)Q_{\text{new}}(x) - P_{\text{new}}(x)| - \delta Q_{\text{new}}(x)}{Q(x)} \right\}.$$

The weighted version of the DC algorithm is discussed, for instance, in [17].

## 2.2 Lattice basis reduction step

Our goal is to compute a simple E-fraction

$$\widehat{R}(x) = \frac{\sum_{j=0}^{\mu} \widehat{p}_j x^j}{1 + \sum_{j=1}^{\nu} \widehat{q}_j x^j},$$

where  $\widehat{p}_j$  and  $\widehat{q}_j$  are fixed-point or floating-point numbers [10, 18], that is as close as possible to  $f_s$ , the function we want to evaluate. These unknown coefficients are of the form  $M2^e$ ,  $M \in \mathbb{Z}$ :

- for fixed-point numbers,  $e$  is implicit (decided at design time);
- for floating-point numbers,  $e$  is explicit (*i.e.*, stored). A floating-point number is of precision  $t$  if  $2^{t-1} \leq M < 2^t - 1$ .

A different format can be used for each coefficient of the desired fraction. If we assume a target format is given for each coefficient, then a straightforward approach is to round each coefficient of  $R_s$  to the desired format. This yields what we call in the sequel a *naive rounding* approximation. Unfortunately, this can lead to a significant loss of accuracy. We first briefly recall the approach from [8] that makes it possible to overcome this issue. Then, we present a small trick that improves on the quality of the output of the latter approach. Eventually, we explain how to handle a coefficient saturation issue appearing in some high radix cases.

### 2.2.1 Modeling with a closest vector problem in a lattice

Every fixed-point number constraint leads to a corresponding unknown  $M$ , whereas each precision- $t$  floating-point number leads to two unknowns  $M$  and  $e$ . A heuristic trick is given in [19] to find an appropriate value for each  $e$  in the floating-point case: we assume that the coefficient in question from  $\widehat{R}$  will have the same order of magnitude as the corresponding one from  $R_s$ , hence they have the same exponent  $e$ . Once  $e$  is set, the problem is reduced to a fixed-point one.

Then, given  $u_0, \dots, u_{\mu}, v_1, \dots, v_{\nu} \in \mathbb{Z}$ , we have to determine  $\mu + \nu + 1$  unknown integers  $a_j (= \widehat{p}_j 2^{-u_j})$  and  $b_j (= \widehat{q}_j 2^{-v_j})$  such that the fraction

$$\widehat{R}(x) = \frac{\sum_{j=0}^{\mu} a_j 2^{u_j} x^j}{1 + \sum_{j=1}^{\nu} b_j 2^{v_j} x^j}$$

is a good approximation to  $R_s$  (and  $f_s$ ), *i.e.*,  $\|\widehat{R} - R_s\|$  is small. To this end, we discretize the latter condition in  $\mu + \nu + 1$  points  $x_0 < \dots < x_{\mu+\nu}$  in the interval  $[a, b]$ , which gives rise to the following instance of a closest vector problem, one of the fundamental questions in the algorithmics of Euclidean lattices [20]: we want to compute  $a_0, \dots, a_{\mu}, b_1, \dots, b_{\nu} \in \mathbb{Z}$  such that the vectors

$$\sum_{j=0}^{\mu} a_j \alpha_j - \sum_{j=1}^{\nu} b_j \beta_j \text{ and } \mathbf{r} \tag{9}$$

are as close as possible, where  $\alpha_j = [2^{u_j} x_0^j, \dots, 2^{u_j} x_{\mu+\nu}^j]^t$ ,  $\beta_j = [2^{v_j} x_0^j R_s(x_0), \dots, 2^{v_j} x_{\mu+\nu}^j R_s(x_{\mu+\nu})]^t$  and  $\mathbf{r} = [R_s(x_0), \dots, R_s(x_{\mu+\nu})]^t$ . It can be solved in an approximate way very efficiently by applying techniques introduced in [21] and [22]. We refer the reader to [8, 19] for more details on this and how the discretization  $x_0, \dots, x_{\mu+\nu}$  should be chosen.

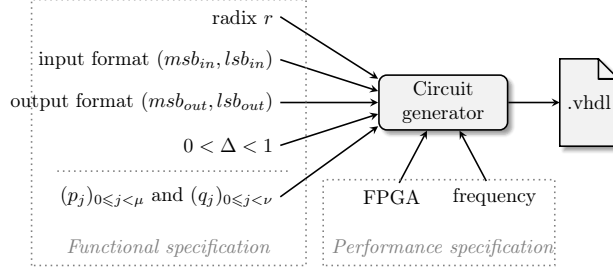


Figure 1: Circuit generator overview.

### 2.2.2 A solution to a coefficient saturation issue

While we generally obtain integer  $a_j$  and  $b_j$  which correspond to a good approximation, the solution is not always guaranteed to give a valid simple E-fraction. What happens is that, in many cases, some of the denominator coefficients in  $R_s$  are maximal with respect to the magnitude constraint in (6) (recall that the second line in (6) can be restated as  $|q_j| \leq \alpha - \max(|a|, |b|)$ ). In this context, the corresponding values of  $|b_j|$  are usually too large. We thus propose to fix the problematic values of  $b_j$  to the closest value to the allowable limit that does not break the E-method magnitude constraints.

The change is minor in (9); we just move the corresponding vectors in the second sum on the left hand side of (9) to the right hand side with opposite sign. The resulting problem can also be solved using the tools from [8, 19]. This usually gives a valid simple E-fraction  $\hat{R}$  of very good quality.

### 2.2.3 Higher radix problems

Coefficient saturation issues get more pronounced by increasing the radix  $r$ . In such cases, care must also be taken with the approximation domain: the  $|q_j|$  upper magnitude bound  $\alpha - \max(|a|, |b|)$  can become negative, since  $\alpha = (1 - \Delta)/(2r) \rightarrow 0$  as  $r \rightarrow \infty$ . To counter this, we use argument and domain scaling ideas presented in [23]. This basically consists in approximating  $f(x) = f(2^t y)$ , for  $y \in [2^{-t}a, 2^{-t}b]$  as a function in  $y$ . If  $t > 0$  is large enough, then the new  $|q_j|$  bound  $\alpha - \max(|2^{-t}a|, |2^{-t}b|)$  will be  $\geq 0$ .

## 3 A hardware implementation targeting FPGAs

We now focus on the hardware implementation of the E-method on FPGAs. This section introduces a generator capable of producing circuits that can solve the system  $\mathbf{A} \cdot \mathbf{y} = \mathbf{b}$ , through the recurrences of Equations (3)–(5).

The popularity of FPGAs is due to their ability to be reconfigured, and their relevance in prototyping as well as in scientific and high-performance computing. They are composed of large numbers of small look-up tables (LUTs), with 4-6 inputs and 1-2 outputs, and each can store the result of any logic function of their inputs. Any two LUTs on the device can communicate, as they are connected through a programmable interconnect network. Results of computations can be stored in registers, usually two of them being connected to the output of each LUT. These features make of FPGAs a good candidate as a platform for implementing the E-method, as motivated even further below.

### 3.1 A minimal interface

An overview of the hardware back-end is presented in Figure 1. Although not typically accessible to the user, its interface is split between the *functional* and the *performance* specification. The former consists of the radix  $r$ , the input and output formats, specified as the weights of their most significant (MSB) and least significant (LSB) bits, the parameter  $\Delta$  (input by the user) and the coefficients of the polynomials  $P_\mu(x)$  and  $Q_\nu(x)$  (coming from the rational approximation). Having  $msb_{in}$  as a parameter is justified by some of the examples of Section 4, where even though the input  $x$  belongs to  $[-1, 1]$ , the maximum value it is allowed to have is smaller, given by the constraints (6) and (7). This allows for optimizing the datapath sizes. The  $msb_{out}$  can be computed during the polynomial/rational approximation phase and passed to the generator.

The circuit generator is developed inside the FloPoCo framework [24], which facilitates the support of classical parameters in the performance specification, such as the target frequency of the circuit, or the target FPGA



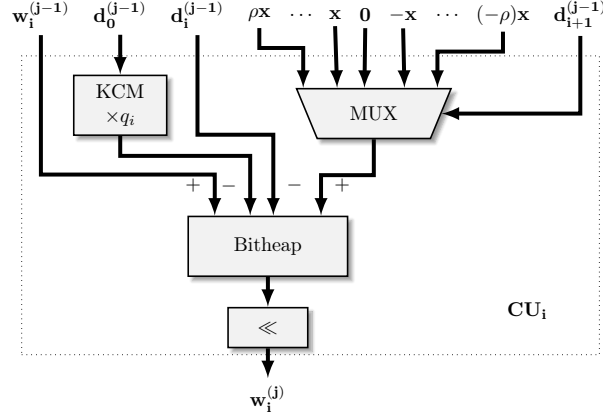


Figure 2: The basic Computation Unit (CU).

device. It also means that we can leverage on the automatic pipelining and test infrastructure present in the framework, alongside the numerous existing arithmetic operators.

### 3.2 Implementation details

An overview of the basic iteration, based on Equation (3), is presented in Figure 2. As this implementation is targeted towards FPGAs, several optimizations can be applied. First, the multiplication  $d_0^{(j-1)} \cdot q_i$  can be computed using the KCM technique for multiplying by a constant [25, 26], with the optimizations of [27], that extend the method for real constants. Therefore, instead of using dedicated multiplier blocks (or of generating partial products using LUTs), we can directly tabulate the result of the multiplication  $d_0^{(j-1)} \cdot q_i$ , at the cost of one LUT per output bit of the result. This remains true even for higher radices, as LUTs on modern devices can accommodate functions of 6 Boolean inputs.

A second optimization relates to the term  $d_{i+1}^{(j-1)} \cdot x$ , from Equation (3). Since  $d_{i+1}^{(j-1)} \in \{-\rho, \dots, \rho\}$ , we can compute the products  $x \cdot \rho$ ,  $x \cdot (\rho - 1)$ ,  $x \cdot (\rho - 2)$ ,  $\dots$ , only once, and then select the relevant one based on the value of  $d_{i+1}^{(j-1)}$ . The multiplications by the negative values in the digit set come at the cost of just one bitwise operation and an addition, which are combined in the same LUT by the synthesis tools.

Finally, regarding the implementation of the CUs, the multi-operand addition of the terms of Equation (3) is implemented using a *bitheap* [28]. The alignments of the accumulated terms and their varied sizes would make for a wasteful use of adders. Using a bitheap we have a single, global optimization of the accumulation. In addition, managing sign extensions comes at the cost of a single additional term in the accumulation, using a technique from [10]. The sign extension of a two's complement fixed point number  $\mathbf{sxx} \dots \mathbf{xx}$  is performed as:

$$\begin{array}{r}
 00 \dots 0 \mathbf{sxx} \dots \mathbf{xxxxxx} \\
 + \quad 11 \dots 110000000 \\
 = \quad \mathbf{ss} \dots \mathbf{sx} \dots \mathbf{xxxxxx}
 \end{array}$$

The sum of the constants is computed in advance and added to the accumulation. The final shift comes at just the cost of some routing, since the shift is by a constant amount.

Modern FPGAs contain fast carry chains. Therefore, we represent the components of the residual vector  $w$  using two's complement, as opposed to a redundant representation. The selection function uses an estimate of  $w_i$  with one fractional bit allowing a simple tabulation using LUTs.

Iteration 0, the initialization, comes at almost no cost, and can be done through fixed routing in the FPGA. This is also true for the second iteration, as simplifying the equations results in  $w_i^{(1)} = r \times w_i^{(0)}$ . The corresponding digits  $d_i^{(1)}$  can be pre-computed and stored directly. This not only saves one iteration, but also improves the accuracy, as  $w_i^{(1)}$  and  $d_i^{(1)}$  can be pre-computed using higher-precision calculations. Going one iteration further, we can see that most of the computations required for  $w_i^{(2)}$  can also be done in advance, except, of course, those involving  $x$ .

Figure 3 shows an unrolled implementation of the E-method, that uses the CUs of Figure 2 as basic building blocks. At the top of the architecture, a scaling can be applied to the input. This step is optional, and the



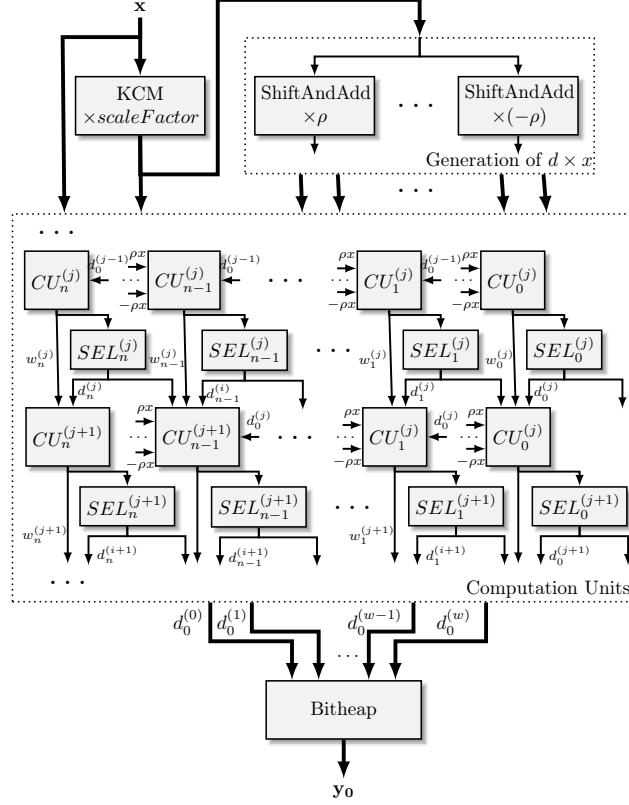


Figure 3: The E-method circuit generator.

scale factor (optional parameter in the design) can either be set by the user, or computed by the generator so that, given the input format, the parameter  $\Delta$  and the coefficients of  $P$  and  $Q$ , the scaled input satisfies the constraints (6) and (7). The multiplications between  $x$  and the possible values of the digits  $d_i^{(j)}$  are done using the classical shift-and-add technique, a choice justified by the small values of the constants and the small number of bits equal to 1 in their representations. At the bottom of Figure 3, the final result  $y_0$  is obtained in two's complement representation. Again, this step is also optional, as users might be content with having the result in the redundant representation.

There is one more optimization that can be done here due to an unrolled implementation. Because only the  $d_0^{(j)}$  digits are required to compute  $y_0$ , after iteration  $m - n$  we can compute one less element of  $\mathbf{w}^{(j)}$  and  $\mathbf{d}^{(j)}$  at each iteration. This optimization is the most effective when the number of required iterations  $m$  is comparable to  $n$ , in which case the required hardware is reduced to almost half.

### 3.3 Error Analysis

To obtain a minimal implementation for the circuit described in Figure 3, we need to size the datapaths in a manner that guarantees that the output  $y_0$  remains unchanged, with respect to an ideal implementation, free of potential rounding errors. To that end, we give an error-analysis, which follows [6, Ch. 2.8]. For the sake of brevity, we focus on the radix 2 case.

In order for the circuit to produce correct results, we must ensure that the rounding errors do not influence the selection function:  $S(\tilde{w}_i^{(j)}) = S(w_i^{(j)}) = d_i^{(j)}$ , where the tilded terms represent approximate values. In [6], the idea is to model the rounding errors due to the limited precision used to store the coefficients  $p_j$  and  $q_j$  inside the matrix  $A$  as a new error matrix  $\mathbf{E}_A = (\varepsilon_{ij})_{n \times n}$ . With the method introduced in this paper, the coefficients are machine representable numbers, and therefore incur no additional error. What remains to deal with are errors due to the limited precision of the involved operators. The only one that could produce rounding errors is the multiplication  $d_0^{(j-1)} \cdot q_i$ . We know that  $d_0^{(j-1)} \geq 1$  (the case  $d_0^{(j-1)} = 0$  is clearly not a problem), so the LSB of  $d_0^{(j-1)} \cdot q_i$  is at least that of  $q_i$ , if not larger. If the output precision satisfies  $lsb_{out} \geq lsb_{q_i}$  (which is usually

Table 1: Approximation errors in the real coefficient and fixed-point coefficient E-fraction cases

	Function Type of error	$\Delta$	$r$	$(\mu, \nu)$	$m$	Real coefficient E-fraction error	Lattice-based error	Naive rounding error
Ex. 1	$\sqrt{1 + (9x/2)^4}, x \in [0, 1/32]$ absolute	$\frac{1}{8}$	2	(4, 4)	32	$5.22 \cdot 10^{-11}$	$5.71 \cdot 10^{-11}$	$1.11 \cdot 10^{-9}$
			4			$6.32 \cdot 10^{-11}$	$4.93 \cdot 10^{-10}$	$7 \cdot 10^{-11}$
			8			$8.25 \cdot 10^{-11}$	$1.11 \cdot 10^{-10}$	$1.78 \cdot 10^{-9}$
Ex. 2	$\exp(x), x \in [0, 7/64]$ relative	$\frac{1}{8}$	2	(3, 3)	32	$1.64 \cdot 10^{-10}$	$1.94 \cdot 10^{-10}$	$3.24 \cdot 10^{-10}$
			4	(4, 4)		$10^{-12}$	$1.11 \cdot 10^{-12}$	$1.91 \cdot 10^{-11}$
			8	(5, 0)		$1.16 \cdot 10^{-12}$	$1.39 \cdot 10^{-12}$	$1.74 \cdot 10^{-11}$
Ex. 3	$\log_2(1 + 2^{-16x}), x \in [0, 1/16]$ absolute	$\frac{1}{2}$	2	(5, 5)	24	$1.98 \cdot 10^{-8}$	$2.33 \cdot 10^{-8}$	$4.37 \cdot 10^{-7}$
			4,8,16	(5, 0)		$2.04 \cdot 10^{-8}$	$2.64 \cdot 10^{-8}$	$4.22 \cdot 10^{-7}$
Ex. 4	$\operatorname{erf}(x), x \in [0, 1/32]$ absolute	$\frac{1}{8}$	2	(4, 4)	48	$2.92 \cdot 10^{-17}$	$3.43 \cdot 10^{-17}$	$1.67 \cdot 10^{-16}$
			4	(4, 4)		$3.44 \cdot 10^{-17}$	$4.23 \cdot 10^{-17}$	$1.13 \cdot 10^{-16}$
			8,16,32	(5, 0)		$1.34 \cdot 10^{-15}$	$1.64 \cdot 10^{-15}$	$2.7 \cdot 10^{-15}$
Ex. 5	$J_0(x), x \in [-1/16, 1/16]$ relative	$\frac{1}{2}$	2	(4, 4)	48	$2.15 \cdot 10^{-17}$	$2.37 \cdot 10^{-15}$	$2.49 \cdot 10^{-15}$
			4,8	(6, 0)		$1.23 \cdot 10^{-17}$	$2.37 \cdot 10^{-15}$	$2.53 \cdot 10^{-15}$

the case), we perform this operation on its full precision, so we do not require any additional guard bits for the internal computations. If this assumption does not hold, based on [6], we obtain the following expression for the rounding errors introduced when computing  $\mathbf{w}^{(j)}$  inside Equation (2), denoted with  $\boldsymbol{\varepsilon}_{\mathbf{w}}^{(j)}$ :

$$\boldsymbol{\varepsilon}_{\mathbf{w}}^{(j)} = 2 \cdot (\boldsymbol{\varepsilon}_{\mathbf{w}}^{(j-1)} + \boldsymbol{\varepsilon}_{const\_mult} + \mathbf{E}_{\mathbf{A}} \cdot \mathbf{d}^{(j-1)}).$$

We can thus obtain an expression for  $\boldsymbol{\varepsilon}_{\mathbf{w}}^{(m)}$ , the error vector at step  $m$ , where  $m$  is the bitwidth of  $y_0$  and  $\boldsymbol{\varepsilon}_{const\_mult}$  are the errors due to the constant multipliers. Since  $\boldsymbol{\varepsilon}_{\mathbf{w}}^{(0)} = 0$ ,

$$\boldsymbol{\varepsilon}_{w_0}^{(m)} = 2^m \cdot (\boldsymbol{\varepsilon}_{const\_mult} + \|\mathbf{E}_{\mathbf{A}}\| \cdot \sum_{j=1}^m d_0^{(j)} \cdot 2^{-j}),$$

where  $\|\mathbf{E}_{\mathbf{A}}\|$  is the matrix 2-norm. We use a larger intermediary precision for the computations, with  $g$  extra guard bits. Therefore, we can design a constant multiplier for which  $\boldsymbol{\varepsilon}_{const\_mult} \leq \bar{\boldsymbol{\varepsilon}}_{const\_mult} \leq 2^{-m-g}$ . Also,

$$\|\mathbf{E}_{\mathbf{A}}\| \leq \max_i \sum_{j=1}^n |\varepsilon_{ij}| \quad \text{and} \quad \sum_{j=1}^m d_0^{(j)} \cdot 2^{-j} < 1,$$

hence we can deduce that for each  $w_i^{(m)}$  we have  $\varepsilon_{w_i}^{(m)} \leq \bar{\varepsilon}_{w_i}^{(m)} \leq 2^m (2^{-m-g} + n \cdot 2^{-m-g})$ . In order for the method to produce correct results, we need to ensure that  $\bar{\varepsilon}_w^{(m)} \leq \Delta/2$ , therefore we need to use  $g \geq 2 + \log_2(2(n+1)/\Delta)$  additional guard bits. This also takes into consideration the final rounding to the output format.

## 4 Examples, Implementation and Discussion

In this section, we consider fractions with fixed-point coefficients of 24, 32, and 48 bits: these coefficients will be of the form  $i/2^m$ , with  $-2^m \leq i \leq 2^m$ , where  $w = 24, 32, 48$ . The target approximation error in each case is  $2^{-m}$ , i.e.,  $\sim 5.96 \cdot 10^{-8}, 2.33 \cdot 10^{-10}, 3.55 \cdot 10^{-15}$  respectively.

**Examples.** All the examples are defined in the first column of Table 1. When choosing them we considered:

- Functions useful in practical applications. The exponential function (Example 2) is a ubiquitous one. Functions of the form  $\log_2(1 + 2^{\pm kx})$  (as the one of Example 3) are useful when implementing logarithmic number systems. The erf function (Example 4) is useful in probability and statistics, while the Bessel function  $J_0$  (Example 5) has many applications in physics.
- Functions that illustrate the various cases that can occur: polynomials are a better choice (Example 3); rational approximation is better (Examples 1, 2, Example 4 if  $r \leq 8$  and Example 5 if  $r = 2$ ). We also include instances where the approximating E-fractions are very different from the minimax, unconstrained, rational approximations with similar degrees in the numerator and denominator (Examples 1 and 2).

All the examples start with a radix 2 implementation after which higher values of  $r$  are considered. Table 1 displays approximation errors in the real coefficient and fixed-point coefficient E-fraction cases. Notice in particular the lattice-based approximation errors, which are generally much better than the naive rounding ones. We also give some complementary comments.

*Example 1.* The type (4, 4) rational minimax unconstrained approximation error is  $4.59 \cdot 10^{-16}$ , around 5 orders of magnitude smaller than the E-fraction error. A similar difference happens in case of Example 2, where the type (3, 3) unconstrained minimax approximation has error  $2.26 \cdot 10^{-16}$ .

*Example 2.* In this case, we are actually working with a rescaled input and are equivalently approximating  $\exp(2x)$ ,  $x \in [0, 7/128]$ . Also, for  $r = 8$ , the real coefficient E-fraction is the same as the E-polynomial one (the magnitude constraint for the denominator coefficients is 0).

*Example 3.* Starting with  $r = 8$ , we have to scale both the argument  $x$  and the approximation domain by suitable powers of 2 for the E-method constraints to continue to hold (see end of Section 2.1).

*Example 4.* As with the previous example, for  $r = 16, 32$  we have to rescale the argument and interval to get a valid E-polynomial.

*Example 5.* By a change of variable, we are actually working with  $J_0(2x - 1/16)$ ,  $x \in [0, 1/16]$ . If we consider  $r \geq 16$ , the 48 bits used to represent the coefficients were not sufficient to produce an approximation with error below  $2^{-48}$ .

**Implementation.** We have generated the corresponding circuits for each of the examples, and synthesized them. The target platform is a Xilinx Virtex6 device xc6vcx75t-2-ff484, and the toolchain used is ISE 14.7. The resulting circuit descriptions are in an easily readable and portable VHDL. For each of the examples compare against a state of the art polynomial approximation implementation generated by FloPoCo (described in [9]). FloPoCo [24] is an open-source arithmetic core generator and library for FPGAs. It is, to the best of our knowledge, one of the only alternatives capable of producing the functions chose for comparison. Table 2 presents the results.

At the top of Table 2, for Example 1, we show the flexibility of the generator: it can achieve a wide range of latencies and target frequencies. The examples show how the frequency can be scaled from around 100MHz to 300MHz, at the expense of a deeper pipeline and an increased number of registers. The number of registers approximately doubles each time the circuit’s period is reduced by a factor 2. This very predictable behavior should help the end user make an acceptable trade-off in terms of performance to required resources. The frequency cap of 300MHz is not something inherent to the E-method algorithm, neither to the implementation. Instead it comes from pipeline performance issues of the bitheap framework inside the FloPoCo generator. We expect that once this bottleneck is fixed, our implementations will reach much higher target frequencies, without the need for any modifications in the current implementation.

**Discussion.** Examples 1 and 2 illustrate that for functions where classical polynomial approximation techniques, like the one used in FloPoCo, manage to find solutions of a reasonably small degree, the ensuing architectures also manage to be highly efficient. This shows, as implementations produced by FloPoCo (with polynomials of degree 6 in both cases) are twice (if not more) as efficient in terms of resources.

However, this is no longer the case when E-fractions can provide a better approximation. This is reflected by Examples 3 to 5, where we obtain a more efficient solution, by quite a large margin in some cases.

For Example 5, Table 2 does not present any data for the FloPoCo implementation as they do not currently support this type of function.

There are a few remarks to be made regarding the use of a higher radix in the implementations of the E-method. Example 4 is an indication that the overall delay of the architecture reaches a point where it can no longer benefit from increasing the radix. The lines of Table 2 marked with an asterisk were generated with an alternative implementation for the CUs, which uses multipliers for computing the  $d_{i+1}^{(j-1)} \cdot x$  products. This is due to the exponential increase of the size of multiplexers with the increase of the radix, while the equivalent multiplier only increases linearly. Therefore, there is a crossover point from which it is best to use this version of the architecture, usually at radix 8 or 16. Finally, the effects of truncating the last iterations become the most obvious when the maximum degree  $n$  is close to the number of required iterations  $m$  in radix  $r$ . This effect can be observed for Example 3 and 4, where there is a considerable drop in resource consumption between the use of radix 8 and 16, and 16 and 32, respectively.

Table 2: Synthesis results for a Xilinx Virtex6 device

Design	Approach	radix	Resources		Performance cycles@period(ns)
			LUT	reg.	
Ex. 1	Ours	2	7,880	0	1@94.3
			7,966	1,523	11@9.6
			7,299	2,689	17@5.7
			6,786	5,202	36@3.7
		4	4,871	0	1@57.9
			4,768	988	7@12.3
			4,600	1,583	11@6.9
			4,853	3,106	22@3.8
		8	4,210	0	1@44.4
			3,875*	0	1@62.2*
			5,307*	309	5@18.4*
			5,184*	499	8@10.4*
			4,707*	1,027	15@5.8*
	FloPoCo	—	994	0	1@29.5
			1,032	138	7@6.7
Ex. 2	Ours	2	6,820	0	1@88.5
		4	6,356	0	1@68.0
		8	5,042	0	1@39.0
	FloPoCo	—	3,024	0	1@41.1
Ex. 3	Ours	2	2,944	0	1@67.0
		4	2,742	0	1@35.1
		8	2,582	0	1@33.1
		16	2,856	0	1@31.2
			1,565*	0	1@29.0*
	FloPoCo	—	3,622	0	1@55.7
Ex. 4	Ours	2	19,564	0	1@139.6
		4	23,052	0	1@92.5
			21,179*	0	1@131.5*
		8	15,388*	0	1@250.7*
		16	12,878*	0	1@76.9*
		32	3,909*	0	1@86.7*
	FloPoCo	—	20,494	0	1@139.9
Ex. 5	Ours	2	19,423	0	1@368.1
		4	13,642	0	1@70.3
		8	18,653	0	1@58.6
	FloPoCo	—	—	—	—

## 5 Summary and Conclusions

A high throughput system for the evaluation of functions by polynomials or rational functions using simple and uniform hardware is presented. The evaluation is performed using the unfolded version of the E-method, with a latency proportional to the precision. An effective computation of the coefficients of the approximations is given and the best strategies (choice of polynomial vs rational approximation, radix of the iterations) investigated. Designs using a circuit generator for the E-method inside the FloPoCo framework are developed and implemented using FPGAs for five different functions of practical interest, using various radices. As it stands, and as our examples show for FPGA devices, the E-method is generally more efficient as soon as the rational approximation is significantly more efficient than the polynomial one. From a hardware standpoint, the results show it is desirable to use the E-method with high radices, usually at least 8. The method also becomes efficient when we manage to find a balance between the maximum degree  $n$  of the polynomial or E-fraction and the number of iterations required for converging to a correct result, which we can control by varying the radix. An open-source implementation of our approach will soon be available online at <https://github.com/sfilip/emethod>.

## Acknowledgments

This work was partly supported by the FastRelax project of the French Agence Nationale de la Recherche, EPSRC (UK) under grants EP/K034448/1 and EP/P010040/1, the Royal Academy of Engineering and Imagination Technologies.

## A Convergence of the E-fraction differential correction algorithm

The E-fraction approximation problem can be formulated in terms of finding functions of the form

$$R(x) = P(x)/Q(x) := \frac{\sum_{j=0}^{\mu} \alpha_j x^j}{1 + \sum_{j=1}^{\nu} \beta_j x^j}, \quad (10)$$

that best approximate  $f$ , where the denominator coefficients are bounded in magnitude, *i.e.*,

$$\max_{1 \leq k \leq \nu} |\beta_k| \leq d > 0, \quad (11)$$

and  $Q(x) > 0, \forall x \in X$ . We denote this set as  $\mathcal{R}_{\mu,\nu}(X)$ . In this case, the differential correction algorithm is defined as follows: given  $R_k = P_k/Q_k \in \mathcal{R}_{\mu,\nu}(X)$  satisfying (10) and (11), find  $R = P/Q$  that minimizes the expression

$$\max_{x \in X} \left\{ \frac{|f(x)Q(x) - P(x)| - \delta_k Q(x)}{Q_k(x)} \right\},$$

subject to

$$\max_{1 \leq k \leq \nu} |\beta_j| \leq d,$$

where  $\delta_k = \max_{x \in X} |f(x) - R_k(x)|$ . If  $R = P/Q$  is not good enough, continue with  $R_{k+1} = R$ .

The convergence properties of this variation are similar to those of the original differential correction algorithm. Let  $\delta^* = \inf_{R \in \mathcal{R}_{\mu,\nu}(X)} \|f - R\|$ , where  $\|g\| = \max_{x \in X} |g(x)|$ .

**Theorem 1.** *If  $Q_k(x) > 0, \forall x \in X$  and if  $\delta_k \neq \delta$ , then  $\delta_{k+1} < \delta_k$  and  $Q_{k+1}(x) > 0, \forall x \in X$ .*

*Proof.* Since  $\delta_k \neq \delta^*$ , there exists  $\bar{R} = \bar{P}/\bar{Q} \in \mathcal{R}_{\mu,\nu}(X)$  s.t.  $\|f - \bar{R}\| = \bar{\delta} < \delta_k$ .

Let  $\Delta = \min_{x \in X} |\bar{Q}(x)|$  and  $M = 1 + d \max_{x \in X} \sum_{k=1}^{\nu} |x^k|$  so that  $|Q(x)| \leq M$  over  $X$  for any rational function satisfying (10) and (11). By the definition of  $P_{k+1}$  and  $Q_{k+1}$ , we have

$$\begin{aligned}
& \max_{x \in X} \left\{ \frac{|f(x)Q_{k+1}(x) - P_{k+1}(x)| - \delta_k Q_{k+1}(x)}{Q_k(x)} \right\} \\
& \leq \max_{x \in X} \left\{ \frac{|f(x)\overline{Q}(x) - \overline{P}(x)| - \delta_k \overline{Q}(x)}{Q_k(x)} \right\} \\
& \leq \max_{x \in X} \left\{ \left[ \left| f(x) - \frac{\overline{P}(x)}{\overline{Q}(x)} \right| - \delta_k \right] \frac{\overline{Q}(x)}{Q_k(x)} \right\} \\
& \leq \max_{x \in X} \left\{ (\bar{\delta} - \delta_k) \frac{\overline{Q}(x)}{Q_k(x)} \right\} \\
& = -(\delta_k - \bar{\delta}) \min_{x \in X} \frac{\overline{Q}(x)}{Q_k(x)} \\
& \leq -\frac{\Delta}{M}(\delta_k - \bar{\delta}).
\end{aligned}$$

This chain of inequalities tells us that

$$\delta_k Q_{k+1}(x)/Q_k(x) \geq \Delta(\delta_k - \bar{\delta})/M, \quad \forall x \in X \quad (12)$$

resulting in  $Q_{k+1}(x) > 0, \forall x \in X$  and

$$\{|f(x)Q_{k+1}(x) - P_{k+1}(x)| - \delta_k Q_{k+1}(x)\}/Q_k(x) < 0, \quad \forall x \in X,$$

leading to

$$|f(x) - P_{k+1}(x)/Q_{k+1}(x)| < \delta_k, \quad \forall x \in X.$$

The conclusion now follows by the definition of  $\delta_k$ .  $\square$

**Theorem 2.**

$$\delta_k \rightarrow \delta^* \text{ as } k \rightarrow \infty.$$

*Proof.* By the previous theorem, we know that  $\{\delta_k\}$  is a monotonically decreasing sequence bounded below by 0. It is thus convergent. Denote its limit by  $\delta$  and assume that  $\delta > \delta^*$ . Then, there exists  $\overline{R} = \overline{P}/\overline{Q} \in \mathcal{R}_{\mu, \nu}(X)$  s.t.

$$\|f - \overline{R}\| = \bar{\delta} < \delta.$$

Let  $\Delta$  and  $M$  be as in the previous theorem. By the aforementioned chain of inequalities, we get

$$\left| f(x) - \frac{P_{k+1}(x)}{Q_{k+1}(x)} \right| - \delta_k \leq -\frac{\Delta}{M}(\delta_k - \bar{\delta}) \frac{Q_k(x)}{Q_{k+1}(x)}, \quad \forall x \in X$$

and

$$\delta_{k+1} - \delta_k \leq -\frac{\Delta}{M}(\delta_k - \bar{\delta}) \frac{Q_k(z_k)}{Q_{k+1}(z_k)}, \quad (13)$$

where  $z_k$  is an element of  $X$  minimizing  $\frac{Q_k(x)}{Q_{k+1}(x)}$ . Since  $\delta_k$  converges to  $\delta > \bar{\delta}$ , (13) implies that

$$\lim_{k \rightarrow \infty} \frac{Q_k(z_k)}{Q_{k+1}(z_k)} = 0. \quad (14)$$

From (12) we get

$$Q_{k+1}(x) \geq Q_k(x) \frac{\Delta}{M}(\delta_k - \bar{\delta})/\delta_k \geq c Q_k(x), \quad \forall x \in X, \quad (15)$$

provided that  $c \leq \frac{\Delta}{M}(1 - \bar{\delta}/\delta) \leq \frac{\Delta}{M}(\delta_k - \bar{\delta})/\delta_k$ . Taking

$$c = \min \left\{ \frac{1}{2}, \frac{\Delta}{M}(1 - \bar{\delta}/\delta) \right\},$$

it follows from (14) that there exists a positive integer  $k_0$  such that

$$\frac{Q_k(z_k)}{Q_{k+1}(z_k)} \leq c^{|X|}, \quad \forall k \geq k_0. \quad (16)$$

Using (15) and (16), we have

$$\frac{Q_k(z_k)}{Q_{k+1}(z_k)} \prod_{x \in X, x \neq z_k} \frac{Q_k(x)}{Q_{k+1}(x)} \leq c^{|X|} \left(\frac{1}{c}\right)^{|X|-1} = c,$$

and

$$\prod_{x \in X} \frac{Q_{k+1}(x)}{Q_k(x)} \geq \frac{1}{c} \geq 2.$$

It follows that  $\prod_{x \in X} Q_k(x)$  diverges as  $k \rightarrow \infty$ . This contradicts the normalization condition (11).  $\square$

## References

- [1] J. E. Volder, “The CORDIC computing technique,” *IRE Transactions on Electronic Computers*, vol. EC-8, no. 3, pp. 330–334, 1959, reprinted in E. E. Swartzlander, *Computer Arithmetic*, Vol. 1, IEEE Computer Society Press, Los Alamitos, CA, 1990.
- [2] W. F. Wong and E. Goto, “Fast hardware-based algorithms for elementary function computations using rectangular multipliers,” *IEEE Trans. Comput.*, vol. 43, no. 3, pp. 278–294, Mar. 1994.
- [3] D. D. Sarma and D. W. Matula, “Faithful bipartite ROM reciprocal tables,” in *Proceedings of the 12th IEEE Symposium on Computer Arithmetic (ARITH-12)*, Knowles and McAllister, Eds. IEEE Computer Society Press, Los Alamitos, CA, Jun. 1995, pp. 17–28.
- [4] M. J. Schulte and J. E. Stine, “Approximating elementary functions with symmetric bipartite tables,” *IEEE Trans. Comput.*, vol. 48, no. 8, pp. 842–847, Aug. 1999.
- [5] F. de Dinechin and A. Tisserand, “Multipartite table methods,” *IEEE Trans. Comput.*, vol. 54, no. 3, pp. 319–330, Mar. 2005.
- [6] M. D. Ercegovac, “A general method for evaluation of functions and computation in a digital computer,” Ph.D. dissertation, Dept. of Computer Science, University of Illinois, Urbana-Champaign, IL, 1975.
- [7] —, “A general hardware-oriented method for evaluation of functions and computations in a digital computer,” *IEEE Trans. Comput.*, vol. C-26, no. 7, pp. 667–680, 1977.
- [8] N. Brisebarre, S. Chevillard, M. D. Ercegovac, J.-M. Muller, and S. Torres, “An Efficient Method for Evaluating Polynomial and Rational Function Approximations,” in *19th IEEE Conference on Application-specific Systems, Architectures and Processors. ASAP’2008*, Leuven (Belgium), 2–4 July 2008, pp. 233–238.
- [9] F. de Dinechin, “On fixed-point hardware polynomials,” Tech. Rep., Oct. 2015, working paper or preprint. [Online]. Available: <https://hal.inria.fr/hal-01214739>
- [10] M. D. Ercegovac and T. Lang, *Digital Arithmetic*. Morgan Kaufmann Publishers, San Francisco, CA, 2004.
- [11] E. W. Cheney, *Introduction to Approximation Theory*, 2nd ed. AMS Chelsea Publishing. Providence, RI, 1982.
- [12] M. J. D. Powell, *Approximation theory and methods*. Cambridge University Press, 1981.
- [13] R. Reemtsen, “Modifications of the first Remez algorithm,” *SIAM J. Numer. Anal.*, vol. 27, no. 2, pp. 507–518, 1990.
- [14] E. W. Cheney and H. L. Loeb, “Two new algorithms for rational approximation,” *Numer. Math.*, vol. 3, no. 1, pp. 72–75, 1961.
- [15] S. N. Dua and H. L. Loeb, “Further remarks on the differential correction algorithm,” *SIAM J. Numer. Anal.*, vol. 10, no. 1, pp. 123–126, 1973.
- [16] I. Barrodale, M. J. D. Powell, and F. K. Roberts, “The differential correction algorithm for rational  $\ell_\infty$ -approximation,” *SIAM J. Numer. Anal.*, vol. 9, no. 3, pp. 493–504, 1972.
- [17] D. Dudgeon, “Recursive filter design using differential correction,” *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 22, no. 6, pp. 443–448, 1974.
- [18] J.-M. Muller, N. Brisebarre, F. de Dinechin, C.-P. Jeannerod, V. Lefèvre, G. Melquiond, N. Revol, D. Stehlé, and S. Torres, *Handbook of Floating-Point Arithmetic*. Birkhäuser, 2010.
- [19] N. Brisebarre and S. Chevillard, “Efficient polynomial  $L^\infty$  approximations,” in *ARITH ’07: Proceedings of the 18th IEEE Symposium on Computer Arithmetic, Montpellier, France*. IEEE Computer Society, 2007, pp. 169–176.
- [20] P. Q. Nguyen and B. Vallée, Eds., *The LLL Algorithm - Survey and Applications*, ser. Information Security and Cryptography. Springer, 2010. [Online]. Available: <http://dx.doi.org/10.1007/978-3-642-02295-1>
- [21] A. K. Lenstra, H. W. Lenstra, and L. Lovász, “Factoring polynomials with rational coefficients,” *Math. Ann.*, vol. 261, no. 4, pp. 515–534, 1982.



- [22] L. Babai, “On Lovász’ lattice reduction and the nearest lattice point problem,” *Combinatorica*, vol. 6, no. 1, pp. 1–13, 1986.
- [23] N. Brisebarre and J.-M. Muller, “Functions approximable by E-fractions,” in *Proc. 38th Asilomar Conference on Signals, Systems and Computers*, Nov. 2004.
- [24] F. de Dinechin and B. Pasca, “Designing custom arithmetic data paths with FloPoCo,” *IEEE Des. Test Comput.*, vol. 28, no. 4, pp. 18–27, Jul. 2011.
- [25] K. Chapman, “Fast integer multipliers fit in FPGAs (EDN 1993 design idea winner),” *EDN magazine*, no. 10, p. 80, May 1993.
- [26] M. J. Wirthlin, “Constant coefficient multiplication using look-up tables,” *VLSI Signal Processing*, vol. 36, no. 1, pp. 7–15, 2004.
- [27] A. Volkova, M. Istoan, F. De Dinechin, and T. Hilaire, “Towards Hardware IIR Filters Computing Just Right: Direct Form I Case Study,” Tech. Rep., 2017. [Online]. Available: <http://hal.upmc.fr/hal-01561052/>
- [28] N. Brunie, F. de Dinechin, M. Istoan, G. Sergeant, K. Illyes, and B. Popa, “Arithmetic core generation using bit heaps,” in *Proceedings of the 23rd Conference on Field-Programmable Logic and Applications*, Sep. 2013, pp. 1–8.